

### **Data Warehouse**

### Program Studi Sistem Informasi Fakultas Ilmu Komputer dan Rekayasa





## **DATA MINING**

## Apa Data Mining?



- Data mining (pencarian pengetahuan dari data)
  - Mengekstrak secara otomatis pola atau pengetahuan yang <u>menarik</u> (tidak sederhana, tersembunyi, tidak diketahui sebelumnya, berpotensi berguna) dari data dalam jumlah sangat besar.



Data Mining adalah usaha penemuan pengetahuan di intelejensia buatan (bidang *machine learning*) atau analisis statistik dengan mencari atau menemukan aturan-aturan, pola-pola dan struktur dari himpunan data yang besar.

### Mengapa Data Mining: Banjir Data



- Twitter: 8000an tweet per detik → 600 juta tweet per hari.
- Facebook: 30 milyar item (link, status, note, foto dst) per bulan. 500 juta user menghabiskan 700 milyar menit per bulan di situs FB.
- Indomaret: 4500an gerai, asumsikan 3 transaksi per menit = 12 juta transaksi per hari se Indonesia.
- Kartu kredit visa: berlaku di 200 negara. 10 ribu transaksi per detik → 850 juta transaksi per hari.

## Apa Datamining?



• Nama alternatif: Knowledge discovery (mining) in databases (KDD), knowledge extraction, data/pattern analysis, data archeology, data dredging, information harvesting, business intelligence dsb

• Keuntungan bagi organisasi yang menerapkan data mining?

## Keuntungan Datamining



- Perusahaan fokus ke informasi yg <u>berharga</u> di datawarehouse/databasenya.
- Meramalkan masa depan → perusahaan dapat mempersiapkan diri

## ontoh:



Midwest grocery chain menggunakan DM untuk menganalisisi pola pembelian: saat pria membeli roti di hari Kamis dan Sabtu, mereka juga membeli minuman.

Analisis lebih lanjut: pembeli ini belanja di hari kamis dan sabtu, tapi di hari kamis jumlah item lebih sedikit. Kesimpulan yang diambil: pembeli membeli minuman untuk dihabiskan saat weekend.

Tindak lanjut: menjual minuman dengan harga full di hari Kamis dan Sabtu. Mendekatkan posisi roti dan minuman.

# Lanjutan..

Jika Anda mempunyai kartu kredit, sudah pasti Anda bakal sering menerima surat berisi brosur penawaran barang atau jasa. Jika Bank pemberi kartu kredit Anda mempunyai 1.000.000 nasabah, dan mengirimkan sebuah (hanya satu) penawaran dengan biaya pengiriman sebesar Rp. 1.000 per buah maka biaya yang dihabiskan adalah Rp. 1 Milyar!! Jika Bank tersebut mengirimkan penawaran sekali sebulan yang berarti 12x dalam setahun maka anggaran yang dikeluarkan per tahunnya adalah Rp. 12 Milyar!! Dari dana Rp. 12 Milyar yang dikeluarkan, berapa persenkah konsumen yang benarbenar membeli? Mungkin hanya 10 %-nya saja. Secara harfiah, berarti 90% dari dana tersebut terbuang sia-sia.

## Lanjutan..



• Dari contoh kasus di atas merupakan salah satu persoalan yang dapat diatasi oleh data mining dari sekian banyak potensi permasalahan yang ada. Data mining dapat menambang data transaksi belanja kartu kredit untuk melihat manakah pembelipembeli yang memang potensial untuk membeli produk tertentu. Mungkin tidak sampai presisi 10%, tapi bayangkan jika kita dapat menyaring 20% saja, tentunya 80% dana dapat digunakan untuk hal lainnya.

#### Contoh



- □Proses untuk menemukan model atau fungsi yang menjelaskan atau membedakan konsep atau kelas data dengan tujuan untuk dapat memprediksi kelas dari suatu objek yang labelnya tidak diketahui
- **□**Contoh : Mendeteksi Penipuan
- □Tujuan : Memprediksi kasus kecurangan transaksi kartu kredit.
  - Pendekatan:
    - Menggunakan transaksi kartu kredit dan informasi dilihat dari atribut account holder
      - Kapan cutomer melakukan pembelian, Dengan cara apa customer membayar, sebarapa sering customer membayar secara tepat waktu, dll
    - Beri nama/tanda transaksi yang telah dilaksanakan sebagai transaksi yang curang atau yang baik. Ini sebagai atribut klass (the class attribute.)
    - Pelajari model untuk class transaksi
    - Gunakan model ini untuk mendeteksi kecurangan dengan meneliti transaksi kartu kredit pada account.

## Contoh Aplikasi



Bank me-mining transaksi customer untuk mengidentifikasi customer yang kemungkinan besar tertarik terhadap produk baru.

Setelah teknik ini digunakan, terjadi peningkatan 20 kali lipat penurunan biaya dibandingkan dengan cara biasa.

## Contoh Aplikasi

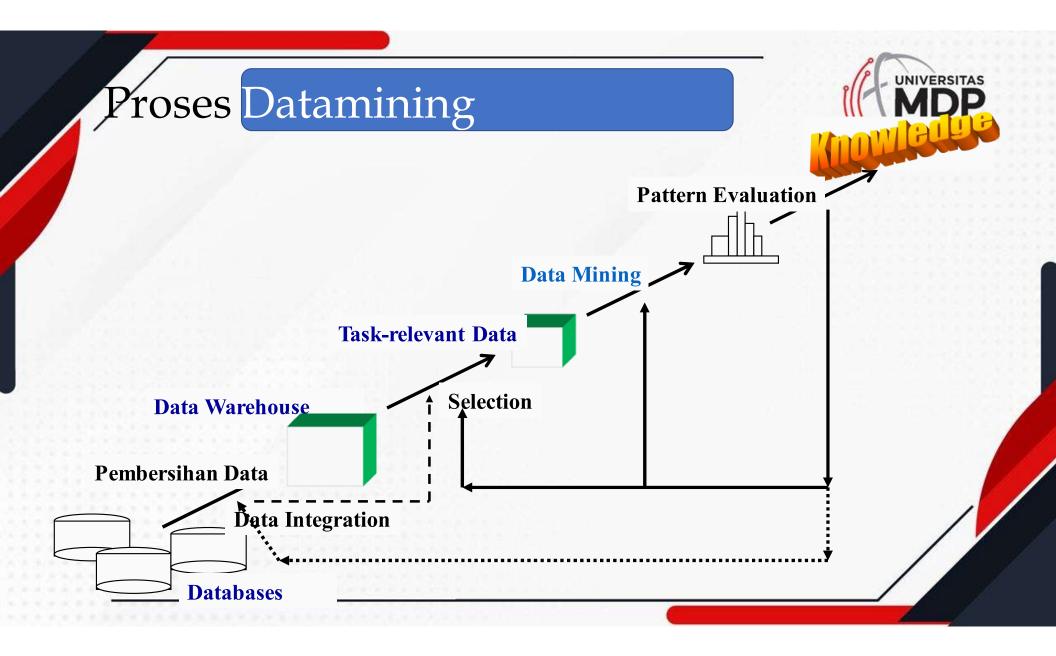


Perusahaan transportasi memining data customer untuk mengelompokkan customer yang memiliki nilai tinggi yang perlu diprioritaskan.



## Data Mining pada Industri Retail

- ■Industri Retail: besarnya data penjualan, sejarah belanja pelanggan, dan lain-lain
- Aplikasi dari Retail data mining
  - ■Mengidentifikasi perilaku pembelian pelanggan
  - ■Menentukan kecenderungan pola belanja pelanggan
  - ■Meningkatkan mutu dari layanan pelanggan
  - ■Mencapai kepuasan pelanggan
  - ■Tingkatkan perbandingan konsumsi barang-barang
  - Mendisain keefektifan distribusi dan transportasi barang





### Data Mining dan Business Intelligence

Semakin mendukung pengambilan keputusan

Pengambilan Keputusan

**Presentasi Data** 

Teknik Visualiasi

**Data Mining** 

Penemuan Informasi

**Eksplorasi Data** 

Statistical Summary, Querying, and Reporting

**Data Preprocessing/Integrasi, Data Warehouses** 

**Sumber Data** 

Database, Web, Paper, Files, Web, eksperimen

**End User** 

**Business Analyst** 

Data Analyst

KR.

## UNIVERSITAS Data Mining: Multi Disiplin Ilmu Teknologi DB Statistik Machine Visualisasi **Data Mining** Learning Pattern Recognition Ilmu Lain Algoritma

### Mengapa tidak analisis data biasa?



- Jumlah data yang sangat besar
  - Algoritma harus scalable untuk menangani data yang sangat besar (tera)
- Dimensi yang sangat besar: ribuan field
- Data Kompleks
  - Aliran data dan sensor
  - Data terstruktur, graph, social network, multi-linked data
  - Database dari berbagai sumber, database lama
  - -Spasial (peta), multimedia, text, web
  - -Software Simulator



### Data Mining dari berbagai sudut pandang



• Relational, data warehouse, web, transactional, stream, OO, spacial, text, multimedia

#### • Pengetahuan yang akan ditambang

• Karakterisitik, diskriminasi, asosiasi, klasifikasi, clustering, trend, outlier

#### Teknik

• Database, OLAP, machine learning, statistik, visualiasi

#### • Penerapan

• Retail, telekomunikasi, banking, analisis kejahatan, bio-data mining, saham, text mining, web mining

### Model dalam Data Mining



#### Verification Model

- Model ini menggunakan (hypothesis) dari pengguna, dan melakukan test terhadap perkiraan yang diambil sebelumnya dengan menggunakan data-data yang ada.
- Model *verifikasi* menggunakan pendekatan *top down* dengan mengambil hipotesa dari user dan memeriksa validitasnya dengan data sehingga bisa dibuktikan kebenaran hipotesa tersebut.

### Model dalam Data Mining



### Discovery Model

- Pada directed knowledge discovery, data mining akan mencoba mencari penjelasan nilai target field tertentu (seperti penghasilan, respons, usia, dan lain-lain) terhadap field-field yang lain.
- Pada *undirected knowledge discovery* tidak ada target field karena komputer akan mencari pola yang ada pada data. Jadi *undirected knowledge discovery* digunakan untuk mengenali hubungan/relasi yang ada pada data sedangkan directed discovery akan menjelaskan hubungan/relasi tersebut.

### Data Mining: Data apa saja?



#### Database Tradisional

- Relational database, data warehouse, transactional database

#### Advanced Database

- Data streams dan data sensor
- Time-series data, temporal data, sequence data (incl. bio-sequences)
- Structure data, graphs, social networks and multi-linked data
- Object-relational databases
- Heterogeneous databases dan legacy databases
- Spatial data dan spatiotemporal data
- Multimedia database
- Text databases
- World-Wide Web

### Top-10 Algorithm di ICDM'06



- #1: C4.5 (61 votes)
- #2: K-Means (60 votes)
- #3: SVM (Support Vector Machine)(58 votes)
- #4: Apriori (52 votes)
- #5: EM (Expectation Maximization) (48 votes)
- #6: PageRank (46 votes)
- #7: AdaBoost (45 votes)
- #7: kNN (45 votes)
- #7: Naive Bayes (45 votes)
- #10: CART (Classification and Regression Tree)(34 votes)

## Aplikasi Data Mining



### Pemasaran/Penyewaan

- Identifikasi pola pembayaran pelanggan
- Menemukan asosiasi diantara karakteristik demografik pelanggan
- Analisis keranjang pemasaran

### Perbankan

- Mendeteksi pola penyalahgunaan kartu kredit
- Identifikasi pelanggan yang loyal
- Mendeteksi kartu kredit yang dihabiskan oleh kelompok pelanggan

### Asuransi & Pelayanan Kesehatan

• Analisis dari klaim

Memprediksi pelanggan yang akan membeli polis baru Identifikasi pola perilaku pelanggan yang berbahaya

## Aplikasi Data Mining



- Analisa Perusahaan dan Manajemen Resiko
  - Perencanaan Keuangan dan Evaluasi Aset
  - Perencanaan Sumber Daya (Resource Planning)
  - Persaingan (competition) → Competitive Intelligence
- Telekomunication
  - menerapkan data mining untuk melihat dari jutaan transaksi yang masuk, transaksi mana saja yang masih harus ditangani secara manual (dilayani oleh orang).

## **Fungsi Data Mining**

- 1. Fungsi Minor atau fungsi tambahan
  - \* Deskription (deskripsi)
  - \* Estimation (estimasi)
  - \* Prediction (prediksi)
- 2. Fungsi Mayor atau fungsi utama
  - \* Classification (klasifikasi)
  - \* Clustering (pengelompokan)
  - \* Association (asosiasi)



### Fungsi Minor



#### Deskripsi

Terkadang peneliti dan analis secara sederhana ingin mencoba mencari cara untuk menggambarkan pola dan kecendrungan yang terdapat dalam data yang dimiliki.

#### • Estimasi

Estimasi hampir sama dengan klasifikasi, kecuali variabel target estimasi lebih ke arah numerik daripada ke arah kategori. Model dibangun menggunakan record lengkap yang menyediakan nilai dari variabel target sebagai nilai prediksi.

#### Prediksi

Prediksi hampir sama dengan klasifikasi dan estimasi, kecuali bahwa dalam prediksi nilal dari hasil akan ada di masa mendatang.

### **Fungsi Mayor**



#### <u>Klasifikasi</u>

Dalam klasifikasi terdapat target variabel kategori, misal penggolongan pendapatan dapat dipisahkan dalam tiga kategori, yaitu tinggi, sedang dan rendah.

#### Pengklusteran

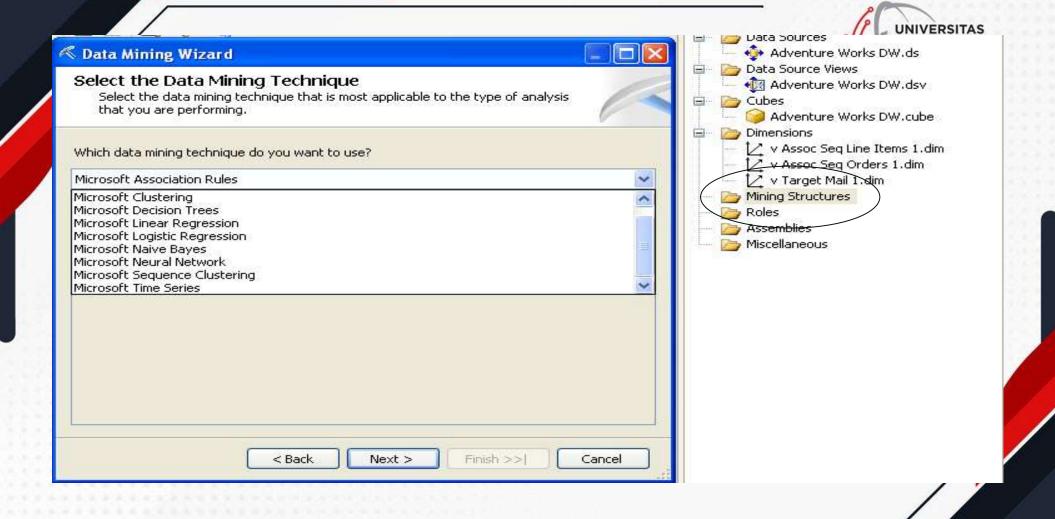
Pengklusteran merupakan pengelompokan record, pengamatan, atau memperhatikan dan membentuk kelas objek-objek yang memiliki kemiripan.

#### Asosiasi

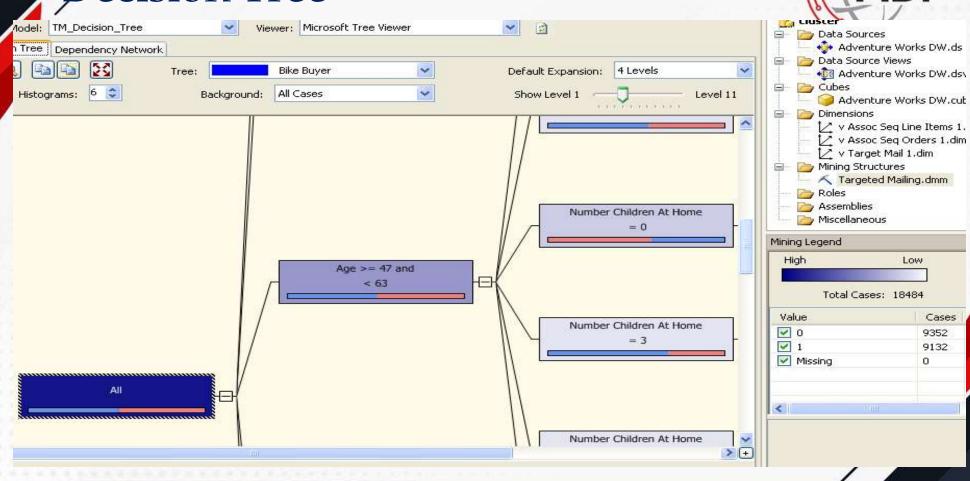
Tugas asosiasi data mining adalah menemukan atribut yang muncul dalam satu waktu. Dalam dunia bisnis lebih umum disebut analisis keranjang belanja



## Data Mining Menggunakan Business Intelligence

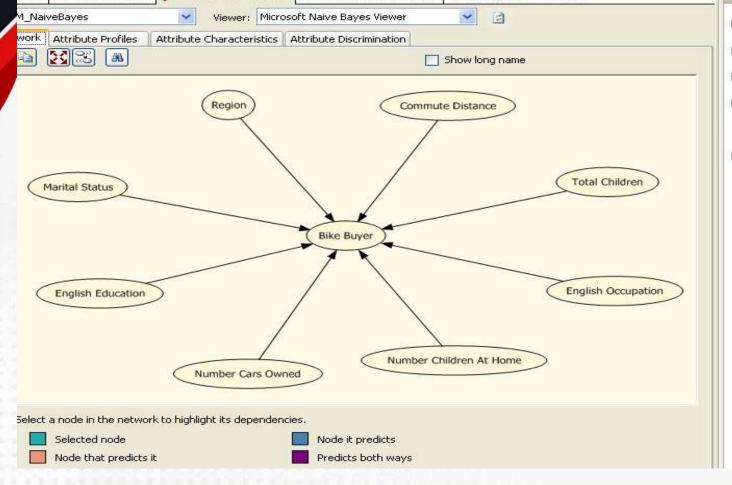




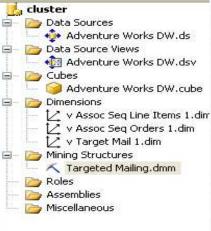


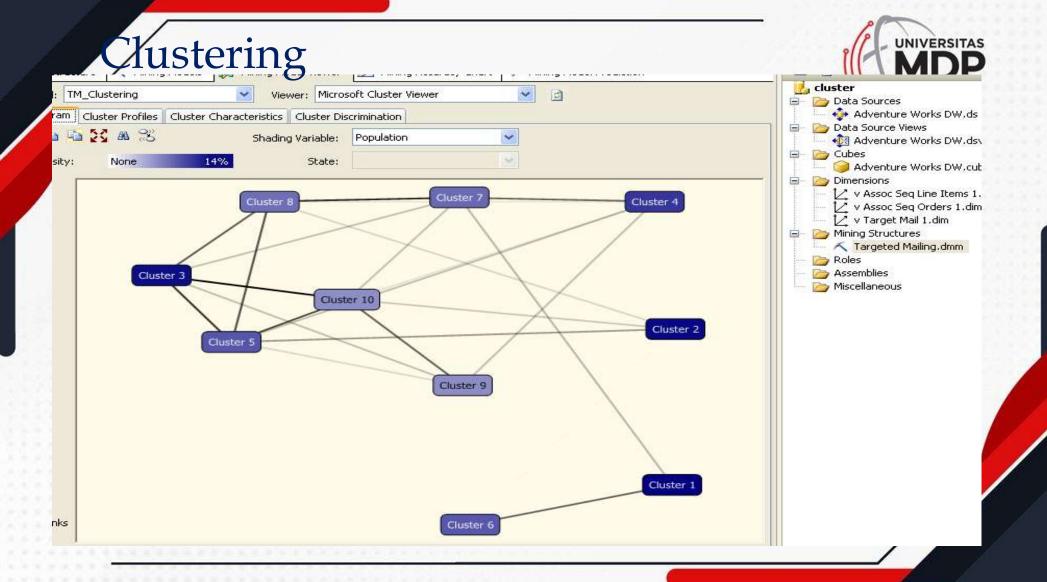
UNIVERSITAS

Vaive Bayes











## **FUNGSI MINOR**

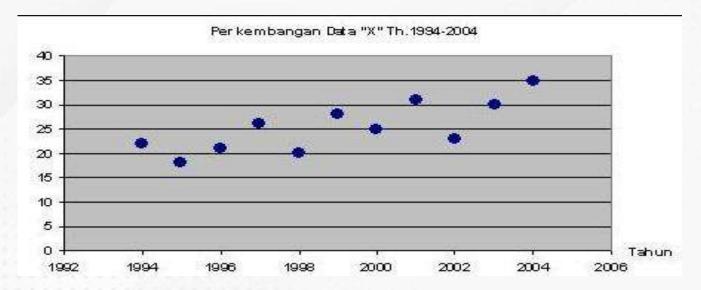
# DescriptionDeskripsi Grafis

- - \* Diagram Titik
  - \* Histogram
- Deskripsi Lokasi
  - \* Rata-rata
  - \* Median
  - \* Modus
  - \* Kuartil, Desil dan Persentil
- Deskripsi Keberagaman
  - Range (rentang)
  - Varians dan Standar Deviasi



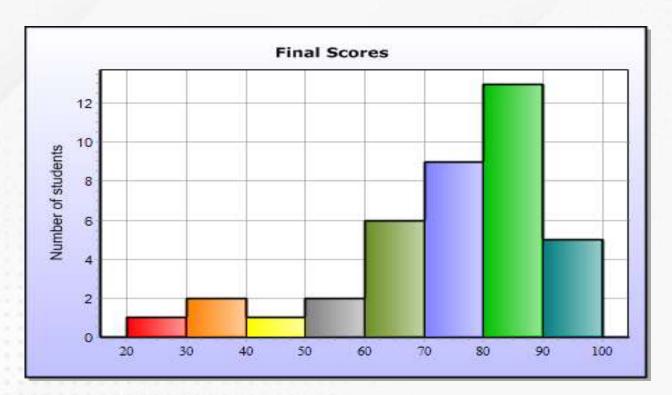
## Diagram Titik





# Histogram





#### Rata-rata



• adalah nilai tunggal yang dianggap dapat mewakili keseluruhan nilai dalam data

$$\mu = \frac{1}{N} \sum_{i=1}^{N} x_i$$

$$\overline{X} = \frac{\sum f_i X_i}{\sum f_i}$$

#### Median



 adalah nilai tengah dari data yang ada setelah data diurutkan

$$Med = X_{\frac{n+1}{2}}$$

$$X_{\frac{n}{2}} + X_{\frac{n}{2}+1}$$

$$Med = L_0 + c \left\{ \frac{\frac{n}{2} - \left(\sum f_1\right)_0}{f_m} \right\}$$

$$Med = \frac{\frac{n}{2} - \left(\sum f_1\right)_0}{f_m}$$

#### Modus



• adalah nilai yang paling sering muncul dalam data

$$Mod = L_0 + c \left\{ \frac{(f_1)_0}{(f_1)_0 + (f_2)_0} \right\}$$

### Kuartil, Desil dan Persentil



• Adalah nilai-nilai yang membagi seperangkat data yang telah terurut menjadi beberapa bagian yang sama

$$Q_{i} = L_{0} + c \left\{ \frac{\frac{in}{4} - \left(\sum f_{i}\right)_{0}}{f_{q}} \right\}$$

$$D_{i} = L_{0} + c \left\{ \frac{\frac{in}{10} - \left(\sum f_{i}\right)_{0}}{f_{d}} \right\}$$

$$P_{i} = L_{0} + c \left\{ \frac{\frac{in}{100} - \left(\sum f_{i}\right)_{0}}{f_{p}} \right\}$$

Range (rentang)



Nilai Jarak = Nilai Maksimum - Nilai Minimum

#### Varians dan Standar Deviasi



$$\sigma^2 = \frac{\sum X^2 - \frac{\left(\sum X\right)^2}{N}}{N}$$

#### 2. Estimation



- Rata-rata sampel sebagai estimasi rata-rata populasi
- Varians sampel sebagai estimasi varians populasi
- Standar Deviasi sampel sebagai standar deviasi populasi

### 3.Prediction



- Regresi Linier Sederhana
- Regresi Linier Berganda

# Régresi Linier Sederhana



$$Y' = a + b X$$

a = Y pintasan, (nilai Y' bila X = 0)

b = kemiringan garis regresi

X = nilai tertentu dari variabel bebas

Y'= nilai yang dihitung pada variabel tidak bebas.

$$b = \frac{\sum_{i} x_{i} y_{i}}{\sum_{i} x_{i}^{2}} \quad a = \overline{Y} - b \overline{X}$$

## Regresi Linier Berganda

$$= b_0 + b_1 X_1 + b_2 X_2 + \dots + b_k X_k$$



$$b_{0}n + b_{1}\sum X_{1} + b_{2}\sum X_{2} + \dots + b_{k}\sum X_{k} = \sum Y$$

$$b_{0}\sum X_{1} + b_{1}\sum X_{1}^{2} + b_{2}\sum X_{1}X_{2} + \dots + b_{k}\sum X_{1}X_{k} = \sum X_{1}Y$$

$$b_{0}\sum X_{2} + b_{1}\sum X_{1}X_{2} + b_{2}\sum X_{2}^{2} + \dots + b_{k}\sum X_{2}X_{k} = \sum X_{2}Y$$

$$\vdots \qquad \vdots \qquad \vdots$$

$$b_0 \sum X_k + b_1 \sum X_1 X_k + b_2 \sum X_2 X_k + \dots + b_k \sum X_k^2 = \sum X_k$$



# TERIMA KASIH

